

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Methods and Systems for Animating Facial Features,
and Methods and Systems for Expression
Transformation**

Inventor(s):

Stephen Marschner

Brian Guenter

Sashi Raghupathy

Kirk Olynyk

Sing Bing Kang

1 **TECHNICAL FIELD**

2 This invention relates to methods and systems for modeling and rendering
3 for realistic facial animation. In particular, the invention concerns methods and
4 systems for facial image processing.

5
6 **BACKGROUND**

7 The field of computer graphics involves rendering various objects so that
8 the objects can be displayed on a computer display for a user. For example,
9 computer games typically involve computer graphics applications that generate
10 and render computer objects for display on a computer monitor or television.
11 Modeling and rendering realistic images is a continuing challenge for those in the
12 computer graphics field. One particularly challenging area within the computer
13 graphics field pertains to the rendering of realistic facial images. As an example, a
14 particular computer graphics application may render a display of an individual
15 engaging in a conversation. Often times, the ultimately rendered image of this
16 individual is very obviously a computer-rendered image that greatly differs from a
17 real individual.

18 Modeling and rendering realistic faces and facial expressions is a
19 particularly difficult task for two primary reasons. First, the human skin has
20 reflectance properties that are not well modeled by the various shading models that
21 are available for use. For example, the well-known Phong model does not model
22 human skin very well. Second, when rendering facial expressions, the slightest
23 deviation from what would be perceived as "real" facial movement is perceived by
24 even the casual viewer as being incorrect. While current facial motion capture
25 systems can be used to create quite convincing facial animation, the captured

motion is much less convincing, and frequently very strange, when applied to another face. For example, if a person provides a sampling of their facial movements, then animating their specific facial movements is not difficult considering that the face from which the movements originated is the same face. Because of this, there will be movement characteristics that are the same or very similar between expressions. Translating this person's facial movements to another person's face, however, is not often times convincing because of, among other things, the inherent differences between the two faces (e.g. size and shape of the face).

Accordingly, this invention arose out of concerns associated with providing improved systems and methods for modeling texture and reflectance of human skin. The invention also arose out of concerns associated with providing systems and methods for reusing facial motion capture data by transforming one person's facial motions into another person's facial motions.

SUMMARY

The illustrated and described embodiments propose inventive techniques for capturing data that describes 3-dimensional (3-D) aspects of a face, transforming facial motion from one individual to another in a realistic manner, and modeling skin reflectance.

In the described embodiment, a human subject is provided and multiple different light sources are utilized to illuminate the subject's face. One of the light sources is a structured light source that projects a pattern onto the subject's face. This structured light source enables one or more cameras to capture data that describes 3-D aspects of the subject's face. Another light source is provided and is

used to illuminate the subject's face. This other light source is sufficient to enable various reflectance properties of the subject's face to be ascertained. The other light source is used in conjunction with polarizing filters so that the specular component of the face's reflectance is eliminated, i.e. only the diffuse component is captured by the camera. The use of the multiple different light sources enables both structure and reflectance properties of a face to be ascertained at the same time. By selecting the light sources carefully, for example, by making the light sources narrowband and using matching narrowband filters on the cameras, the influence of ambient sources of illumination can be eliminated.

Out of the described illumination process, two useful items are produced— (1) a range map (or depth map) and (2) an image of the face that does not have the structured light source pattern in it. A 3D surface is derived from the range map and surface normals to the 3D surface are computed. The processing of the range map to define the 3D surface can optionally include a filtering step in which a generic face template is combined with the range map to reject undesirable noise. The computed surface normals and the image of the face are then used to derive an albedo map. An albedo map is a special type of texture map in which each sample describes the diffuse reflectance of the surface of a face at a particular point on the surface. Accordingly, at this point in the process, information has been ascertained that describes the 3D-aspects of a face (i.e. the surface normals), and information that describes the face's reflectance (i.e. the albedo map).

In one embodiment, the information or data that was produced in the illumination process is used to transform facial expressions of one person into facial expressions of another person. In this embodiment, the notion of a code book is introduced and used.

1 A code book contains data that describes many generic expressions of
2 another person (person A). One goal is to take the code book expressions and use
3 them to transform the expressions of another person (person B). To do this, an
4 inventive method uses person B to make a set of training expressions. The
5 training expressions consist of a set of expressions that are present in the code
6 book. By using the training expressions and each expression's corresponding code
7 book expression, a transformation function is derived. The transformation
8 function is then used to derive a set of synthetic expressions that should match the
9 expressions of person B. That is, once the transformation function is derived, it is
10 applied to each of the expressions in the code book so that the code book
11 expressions match the expressions of person B. Hence, when a new expression is
12 received, e.g. from person B, that might not be in the training set, the synthesized
13 code book expressions can be searched for an expression that best matches the
14 expression of person B.

15 In another embodiment, a common face structure is defined that can be
16 used to transform facial expressions and motion from one face to another. In the
17 described embodiment, the common face structure comprises a coarse mesh
18 structure or "base mesh" that defines a subdivision surface that is used as the basis
19 for transforming the expressions of one person into another. A common base
20 mesh is used for all faces thereby establishing a correspondence between two or
21 more faces. Accordingly, this defines a structure that can be used to adapt face
22 movements from one person to another. According to this embodiment, a
23 technique is used to adapt the subdivision surface to the face model of a subject.
24 The inventive technique involves defining certain points on the subdivision
25 surface that are mapped directly to corresponding points on the face model. This

is true for every possible different face model. By adding this constraint, the base mesh has a property in that it fits different face models in the same way. In addition, the inventive algorithm utilizes a smoothing functional that is minimized to ensure that there is a good correspondence between the base mesh and the face model.

In another embodiment, a reflectance processing technique is provided that gives a measure of the reflectance of the surface of a subject's face. To measure reflectance, the inventive technique separates the reflectance into its diffuse and specular components and focuses on the treatment of the diffuse components.

To measure the diffuse component, an albedo map is first defined. The albedo map is defined by first providing a camera and a subject that is illuminated by multiple different light sources. The light sources are filtered by polarizing filters that, in combination with a polarizing filter placed in front of the camera, suppress specular reflection or prevent specular reflection from being recorded. A sequence of images is taken around the subject's head. Each individual image is processed to provide an individual albedo map that corresponds to that image. All of the albedo maps for a particular subject are then combined to provide a single albedo map for the subject's entire face.

BRIEF DESCRIPTION OF THE DRAWINGS

Fig. 1 is a high level diagram of a general purpose computer that is suitable for use in implementing the described embodiments.

Fig. 2 is a schematic diagram of a system that can be utilized to capture both structural information and reflectance information of a subject's face at the same time.

1 Fig. 15 is a color diagram of the Fig. 14 albedo maps after editing.

2 Fig. 16 is a collection of color pictures of a face model that is rendered in
3 different orientations and under different lighting conditions.

4 Fig. 17 is a flow diagram that describes steps in a method for creating an
5 albedo map in accordance with the described embodiment.

6 Fig. 18 is a flow diagram that describes steps in a method for computing an
7 albedo for a single pixel in accordance with the described embodiment.

8 9 **DETAILED DESCRIPTION**

10 **Overview**

11 Rendering realistic faces and facial expressions requires very good models
12 for the reflectance of skin and the motion of the face. Described below are
13 methods and techniques for modeling, animating, and rendering a face using
14 measured data for geometry, motion, and reflectance that realistically reproduces
15 the appearance of a particular person's face and facial expressions. Because a
16 complete model is built that includes geometry and bi-directional reflectance, the
17 face can be rendered under any illumination and viewing conditions. The
18 described modeling systems and methods create structured face models with
19 correspondences across different faces, which provide a foundation for a variety of
20 facial animation operations.

21 The inventive embodiments discussed below touch upon each of the parts
22 of the face modeling process. To create a structured, consistent representation of
23 geometry that forms the basis for a face model and that provides a foundation for
24 many further face modeling and rendering operations, inventive aspects extend
25 previous surface fitting techniques to allow a generic face to be conformed to

different individual faces. To create a realistic reflectance model, the first known practical use of recent skin reflectance measurements is made. In addition, newly measured diffuse texture maps have been added using an improved texture capture process. To animate a generic mesh, improved techniques are used to produce surface shapes suitable for high quality rendering.

Exemplary Computer System

Preliminarily, Fig. 1 shows a general example of a desktop computer 130 that can be used in accordance with the described embodiments. Various numbers of computers such as that shown can be used in the context of a distributed computing environment. These computers can be used to render graphics and process images in accordance with the description given below.

Computer 130 includes one or more processors or processing units 132, a system memory 134, and a bus 136 that couples various system components including the system memory 134 to processors 132. The bus 136 represents one or more of any of several types of bus structures, including a memory bus or memory controller, a peripheral bus, an accelerated graphics port, and a processor or local bus using any of a variety of bus architectures. The system memory 134 includes read only memory (ROM) 138 and random access memory (RAM) 140. A basic input/output system (BIOS) 142, containing the basic routines that help to transfer information between elements within computer 130, such as during start-up, is stored in ROM 138.

Computer 130 further includes a hard disk drive 144 for reading from and writing to a hard disk (not shown), a magnetic disk drive 146 for reading from and writing to a removable magnetic disk 148, and an optical disk drive 150 for

reading from or writing to a removable optical disk 152 such as a CD ROM or other optical media. The hard disk drive 144, magnetic disk drive 146, and optical disk drive 150 are connected to the bus 136 by an SCSI interface 154 or some other appropriate peripheral interface. The drives and their associated computer-readable media provide nonvolatile storage of computer-readable instructions, data structures, program modules and other data for computer 130. Although the exemplary environment described herein employs a hard disk, a removable magnetic disk 148 and a removable optical disk 152, it should be appreciated by those skilled in the art that other types of computer-readable media which can store data that is accessible by a computer, such as magnetic cassettes, flash memory cards, digital video disks, random access memories (RAMs), read only memories (ROMs), and the like, may also be used in the exemplary operating environment.

A number of program modules may be stored on the hard disk 144, magnetic disk 148, optical disk 152, ROM 138, or RAM 140, including an operating system 158, one or more application programs 160, other program modules 162, and program data 164. A user may enter commands and information into computer 130 through input devices such as a keyboard 166 and a pointing device 168. Other input devices (not shown) may include a microphone, joystick, game pad, satellite dish, scanner, and one or more cameras, or the like. These and other input devices are connected to the processing unit 132 through an interface 170 that is coupled to the bus 136. A monitor 172 or other type of display device is also connected to the bus 136 via an interface, such as a video adapter 174. In addition to the monitor, personal computers typically include other peripheral output devices (not shown) such as speakers and printers.

below can be utilized in various other areas. For example, areas of application include, without limitation, recognition of faces for security, personal user interaction, etc., building realistic face models for animation in games, movies, etc., and allowing a user to easily capture his/her own face for use in interactive entertainment or business communication.

Fig. 2 shows an exemplary system 200 that is suitable for use in simultaneously or contemporaneously capturing facial structure and reflectance properties of a subject's face. The system includes a data-capturing system in the form of one or more cameras, an exemplary one of which is camera 202. Camera 202 can include a CCD image sensor and related circuitry for operating the array, reading images from it, converting the images to digital form, and communicating those images to the computer. The system also includes a facial illumination system in the form of multiple light sources or projectors. In the case where multiple cameras are used, they are genlocked to allow simultaneous capture in time. In the illustrated example, two light sources 204, 206 are utilized. Light source 204 desirably produces a structured pattern that is projected onto the subject's face. Light source 204 can be positioned at any suitable location. This pattern enables structural information or data pertaining to the 3-D shape of the subject's face to be captured by camera 202. Any suitable light source can be used, although a pattern composed of light in the infrared region can be advantageously employed. Light source 206 desirably produces light that enables camera 202 to capture the diffuse component of the face's reflectance property. Light source 206 can be positioned at any suitable location although it has been advantageously placed in line with the camera's lens 202a through, for example, beam splitting techniques. This light source could also be adapted so that it

encircles the camera lens. This light source is selected so that the specular component of the reflectance is suppressed or eliminated. In the illustrated example, a linear polarizing filter is employed to produce polarized illumination, and a second linear polarizer, which is oriented perpendicularly to the first, is placed in front of the lens 202a so that specular reflection from the face is not recorded by the camera. The above-described illumination system has been simulated using light sources at different frequencies, e.g. corresponding to the red and green channels of the camera. Both of the channels can, however, be in the infrared region. Additionally, by selecting the light sources to be in a narrow band (e.g. 780-880 nm), the influence of ambient light can be eliminated. This property is only achieved when the camera is also filtered to a narrow band. Because the illumination from the light source is concentrated into a narrow band of wavelengths whereas the ambient light is spread over a broad range of wavelengths, the light from the source will overpower the ambient light for those particular wavelengths. The camera, which is filtered to record only the wavelengths emitted by the source, will therefore be relatively unaffected by the ambient light. As a result, the camera will only detect the influence of the selected light sources on the subject.

Using the multiple different light sources, and in particular, an infrared light source in combination with a polarized light source (which can be an infrared light source as well) enables the camera (which is configured with a complementary polarizer) to simultaneously or contemporaneously capture structural information or data about the face (from light source 204) and reflectance information or data about the face (from light source 206) independently. The structural information describes 3-dimensional aspects of the face while the reflectance information

Zippered Polygon Meshes from Range Images, SIGGRAPH 94; F. Bernardini, J. Mittleman, H. Rushmeier, C. Silva, and G. Taubin, *The Ball-Pivoting Algorithm for Surface Reconstruction*, Trans. Vis. Comp. Graph. 5:4 (1999). Step 308 then computes surface normal vectors (“surface normals”) to the 3D surface of step 306 using known algorithms. One way to accomplish this task is to compute the normals to the triangles, average those triangle normals around each vertex to make vertex normals, and then interpolate the vertex normals across the interior of each triangle. Other methods can, of course, be utilized. Step 310 then uses the computed surface normals of step 308 and the image data of step 302 to derive an albedo map. An albedo is a special type of texture map in which each sample describes the diffuse reflectance of the surface of a face at a particular point on the facial surface. The derivation of an albedo map, given the information provided above, will be understood by those skilled in the art. An exemplary algorithm is described in Marschner, *Inverse Rendering for Computer Graphics*, PhD thesis, Cornell University, August 1998.

At this point, and as shown in Fig. 2, the illumination processing has produced 3D data that describes the structural features of a subject’s face and albedo map data that describes the diffuse reflectance of the facial surface.

The above illumination processing can be used to extract the described information, which can then be used for any suitable purpose. In one particularly advantageous embodiment, the extracted information is utilized to extract and recognize a subject's expressions. This information can then be used for expression transformation. In the inventive embodiment described just below, the expressions of one person can be used to transform the expressions of another person in a realistic manner.

Expression Transformation Using a Code Book

In one expression transformation embodiment, the notion of a code book is introduced and is utilized in the expression transformation operation that is described below. Fig. 4 shows an exemplary code book 400 that contains many different expressions that have been captured from a person. These expressions can be considered as generic expressions, or expressions from a generic person rather than from a particular individual. In the example, the expressions range from Expression 1 through Expression N. Expression 1 could be, for example, a smile; Expression 2 could be a frown; Expression 3 could be an “angry” expression, and the like. The expressions that are contained in code book 400 are mathematically described in terms of their geometry and can be captured in any suitable way such as the process described directly above.

To effect expression transformation, a transformation function is first derived using some of the expressions in code book 400. To derive the transformation function, the notion of a training set of expressions 402 is introduced. The expression training set 402 consists of a set of expressions that are provided by an individual other than the individual whose expressions are described in the code book 400. The training expressions of training set 402 are a subset of the code book expressions. That is, each expression in the training set corresponds to an expression in the code book 400. For example, the training set 402 might consist of three expressions—Expression 1, Expression 2, and Expression 3, where the expressions are “smile”, “frown” and “angry” respectively. The goal of the transformation function is to take the geometric deformations that are associated with expressions of the training set, and apply

algorithm) will not precisely correspond because of errors in placement. Second, head shape and size varies from person to person.

The first mismatch can be overcome by resampling the motion capture displacement data for all faces at a fixed set of positions on a generic mesh. This is described below in more detail in the section entitled "Exemplary System and Method for Building a Face Model." There, the fixed set of positions is referred to as the "standard sample positions". The resampling function is the mesh deformation function. The standard sample positions are the vertices of the face mesh that correspond to the vertices of the generic mesh subdivided once.

The second mismatch requires transforming displacement data from one face to another to compensate for changes in size and shape of the face. In the illustrated example, this is done by finding a small training set of corresponding expressions for the two data sets and then finding the best linear transformation from one to another. As an example, consider the following: In an experimental environment, emotion expressions were manually labeled for 49 corresponding expressions including various intensities of several expressions. For speech motion, 10,000 frames were automatically aligned using time warping techniques.

Each expression is represented by a $3m$ -vector g that contains all of the x , y , and z displacements at the m standard sample positions. Given a set of n expression vectors for the face to be transformed, $g_{a1...n}$, and a corresponding set of vectors for the target face, $g_{b1...n}$, a set of linear predictors a_j is computed, one for each coordinate of g_a , by solving $3m$ linear least squares systems:

$$a_j \cdot g_{ai} = g_{bi}[j] \quad i = 1...n$$

Exemplary Application

Fig. 6 shows a system 600 that illustrates but one example of how the expression transformation process described above can be employed. System 600 includes a transmitter computing system or transmitter 602 and a receiver computing system or receiver 604 connected for communication by a network 603 such as the Internet. Transmitter 602 includes an illumination system 200 (Fig. 2) that is configured to capture the expressions of a person as described in connection with Fig. 2. Transmitter 602 also includes a code book 400, such as the one described in connection with Fig. 4. It is assumed that the code book has been synthesized into a synthetic set of expressions as described above. That is, using a training set of expressions provided by the person whose expressions illumination system 200 is configured to capture, the code book has been processed to provide the synthesized set of expressions.

Receiver 604 includes a reconstruction module 606 that is configured to reconstruct facial images from data that is received from transmitter 602. Receiver 604 also includes a code book 400 that is identical to the code book that is included with the transmitter 602. Assume now, that the person located at transmitter 602 attempts to communicate with a person located at receiver 604. As the person located at the transmitter 602 moves their face to communicate, their facial expressions and movement are captured and processed by the transmitter 602. This processing can include capturing their expressions and searching the synthesized code book to find the nearest matching expression in the code book. When a matching expression is found in the synthesized code book, an index of

that expression can be transmitted to receiver 604 and an animated face can be reconstructed using the reconstruction module 606.

Exemplary Facial Transformation

Fig. 7 shows some effects of expression transfer in accordance with the described embodiment. The pictures in the first row constitute a synthetic face of a first person (person A) that shows three different expressions. These pictures are the result of the captured facial motion of person A. Face motion for a second person (person B) was captured. The captured face motion for person B is shown in the third row. Here, the 3D motion data was captured by placing a number of colored dots on the person's face and measuring the dots' movements when the person's face was deformed, as will be understood by those of skill in the art. Motion data can, however, be captured by the systems and methods described above. Person B's captured motions were then used, as described above, to transform the expressions of person A. The result of this operation is shown in the second row. The expressions in the three sets of pictures all correspond with one another. Notice how the expressions in the first and second row look very similar even though they were derived from two very different people, while the original expressions of the second person (row 3) look totally unlike those of the first and second rows.

Exemplary System and Methods for Building a Face Model

The model of a face that is needed to produce a realistic image has two parts to it. The first part of the model relates to the geometry of the face (i.e. the shape of the surface of the face) while the second part of the model relates to the

1 reflectance of the face (i.e. the color and reflective properties of the face). This
2 section deals with the first part of that model—the geometry of the face.

3 The geometry of the face consists of a skin surface plus additional surfaces
4 for the eyes. In the present example, the skin surface is derived from a laser range
5 scan of the head and is represented by a subdivision surface with displacement
6 maps. The eyes are a separate model that is aligned and merged with the skin
7 surface to produce a complete face model suitable for high quality rendering.

8 9 **Mesh Fitting**

10 The first step in building a face model is to create a subdivision surface that
11 closely approximates the geometry measured by the laser range scanner. In the
12 illustrated example, the subdivision surfaces are defined from a coarse triangle
13 mesh using Loop's subdivision rules. Loop's subdivision rules are described in
14 detail in Charles Loop, *Smooth Subdivision Surfaces Based on Triangles*, PhD
15 thesis, University of Utah, August 1987. In addition, the subdivision surfaces are
16 defined with the addition of sharp edges similar to those described by Hoppe et al.,
17 *Piecewise smooth surface reconstruction*, Computer Graphics (SIGGRAPH '94
18 Proceedings) pps. 295-302, July 1994. Note that the non-regular crease masks are
19 not used. In addition, when subdividing an edge between a dart and a crease
20 vertex, only the new edge adjacent the crease vertex is marked as a sharp edge.

21 A single base mesh is used to define the subdivision surfaces for all of the
22 face models, with only the vertex positions varying to adapt to the shape of each
23 different face. In the illustrated example, a base mesh having 227 vertices and 416
24 triangles was defined to have the general shape of a face and to provide greater
25 detail near the eyes and lips, where the most complex geometry and motion occur.

The mouth opening is a boundary of the mesh, and is kept closed during the fitting process by tying together the positions of the corresponding vertices on the upper and lower lips. The base mesh has a few edges marked for sharp subdivision rules that serve to create corners at the two sides of the mouth opening and to provide a place for the sides of the nose to fold. Because the modified subdivision rules only introduce creases for chains of at least three sharp edges, this model does not have creases in the surface; only isolated vertices fail to have well-defined limit normals.

Fig. 8 shows an example of a coarse defined mesh (the center figure) that was used in accordance with this example. Fig. 8 visually shows how the coarse mesh can be used to map the same subdivision control (coarse) mesh to a displaced subdivision surface for each face so that the result is a natural correspondence from one face to another. This aspect is discussed in more detail below.

The process used to fit the subdivision surface to each face is based on an algorithm described by Hoppe et al. *Piecewise smooth surface reconstruction*, Computer Graphics (SIGGRAPH '94 Proceedings) pps. 295-302, July 1994. Hoppe's surface fitting method can essentially be described as consisting of three phases: a topological type estimation (phase 1), a mesh optimization (phase 2), and a piecewise smooth surface optimization (phase 3).

Phase 1 constructs a triangular mesh consisting of a relatively large number of triangles given an unorganized set of points on or near some unknown surface. This phase determines the topological type of the surface and produces an initial estimate of geometry. Phase 2 starts with the output of phase 1 and reduces the number of triangles and improves the fit to the data. The approach is to cast the

problem as optimization of an energy function that explicitly models the trade-off between the competing goals of concise representation and good fit. The free variables in the optimization procedure are the number of vertices in the mesh, their connectivity, and their positions. Phase 3 starts with the optimized mesh (a piecewise linear surface) that is produced in phase 2 and fits an accurate, concise piecewise smooth subdivision surface, again by optimizing an energy function that trades off conciseness and fit to the data. The phase 3 optimization varies the number of vertices in the control mesh, their connectivity, their positions, and the number and locations of sharp features. The automatic detection and recovery of sharp features in the surface is an essential part of this phase.

In the present embodiment, processing differs from the approach described in Hoppe et al. in a couple of ways. First, continuous optimization is performed only over vertex positions, since we do not want to alter the connectivity of the control mesh. Additionally, feature constraints are added as well as a smoothing term.

In the illustrated example, the fitting process minimizes the functional:

$$E(\mathbf{v}) = E_d(\mathbf{v}, \mathbf{p}) + \lambda E_s(\mathbf{v}) + \mu E_c(\mathbf{v})$$

where \mathbf{v} is a vector of all the vertex positions, and \mathbf{p} is a vector of all the data points from the range scanner. The subscripts on the three terms stand for distance, shape, and constraints. The distance functional E_d measures the sum-squared distance from the range scanner points to the subdivision surface:

$$E_d(\mathbf{v}, \mathbf{p}) = \sum_{i=1}^{n_p} a_i \|p_i - \Pi(\mathbf{v}, p_i)\|^2$$

$$a_i = \begin{cases} 1 & \text{if } \langle s(p_i), n(\Pi(\mathbf{v}, p_i)) \rangle > 0 \text{ and } \|p_i - \Pi(\mathbf{v}, p_i)\| < d_0 \\ 0 & \text{otherwise} \end{cases}$$

The smoothness functional E_s encourages the control mesh to be locally planar. It measures the distance from each vertex to the average of the neighboring vertices:

$$E_s(\mathbf{v}) = \sum_{j=1}^{n_v} \left\| v_j - \frac{1}{\deg(v_j)} \sum_{i=1}^{\deg(v_j)} v_{k_i} \right\|^2$$

The constraint functional E_c is simply the sum-squared distance from a set of constrained vertices to a set of corresponding target positions:

$$E_c(\mathbf{v}) = \sum_{i=1}^{n_c} \|A_{ci}\mathbf{v} - d_i\|^2$$

0809001521 MSI-546US.PAT.APP.DOC

mapped directly to corresponding points that are marked on the face model. The mapping of these specific points takes place in the same manner for each of the many different possible face models. Step 1006 fits the generic face model to the one or more face models. This step is implemented by manipulating only the positions of the vertices to adapt to the shape of each different face. During the fitting process continuous optimization is performed only over the vertex positions so that the connectivity of the mesh is not altered. In addition, the fitting process involves mapping the specific points or constraints directly to the face model. In addition, a smoothing term is added and minimized so that the control mesh is encouraged to be locally planar.

Adding Eyes

The displaced subdivision surface just described represents the shape of the facial skin surface quite well. There are, however, several other features that are desirable for a realistic face. The most important of these is the eyes. Since the laser range scanner does not capture suitable information about the eyes, the mesh is augmented for rendering by adding separately modeled eyes. Unlike the rest of the face model, the eyes and their motions are not measured from a specific person, so they do not necessarily reproduce the appearance of the real eyes. However, their presence and motion is critical to the overall appearance of the face model.

Any suitable eye model can be used to model the eyes. In the illustrated example, a commercial modeling package was used to build a model consisting of two parts. The first part is a model of the eyeball, and the second part is a model of the skin surface around the eye, including the eyelids, orbit, and a portion of the

Moving the Face

The motions of the face are specified by the time-varying 3D positions of a set of sample points on the face surface. When the face is controlled by motion-capture data these points are the markers on the face that are tracked by the motion capture system. The motions of these points are used to control the face surface by way of a set of control points that smoothly influence regions of the surface. Capturing facial motion data can be done in any suitable way, as will be apparent to those of skill in the art. In one specific example, facial motion was captured using the technique described in Guenter et al., *Making Faces*, Proceedings of SIGGRAPH 1998, pages 55-67, 1998.

Mesh Deformation

The face is animated by displacing each vertex w_i of the triangle mesh from its rest position according to a linear combination of the displacements of a set of control points q_j . These control points correspond one-to-one with the sample points p_j that describe the motion. The influence of each control point on the vertices falls off with distance from the corresponding sample point, and where multiple control points influence a vertex, their weights are normalized to sum to 1.

$$\Delta w_i = \frac{1}{\beta_i} \sum_j \alpha_{ij} \Delta q_j \quad ; \alpha_{ij} = h(\|w_i - p_j\|/r)$$

where $\beta_i = \sum_k \alpha_{ik}$ if vertex i is influenced by multiple control points and 1 otherwise. These weights are computed once, using the rest positions of the sample points and face mesh, so that moving the mesh for each frame is just a

1 sparse matrix multiplication. For the weighting function, the following was used:

2
$$h(x) = \frac{1}{2} + \frac{1}{2}\cos(\pi x).$$

3 Two types of exceptions to these weighting rules are made to handle the
4 particulars of animating a face. Vertices and control points near the eyes and
5 mouth are tagged as "above" and "below," and control points that are, for example,
6 above the mouth do not influence the motions of vertices below the mouth. Also,
7 a scalar texture map in the region around the eyes is used to weight the motions so
8 that they taper smoothly to zero at the eyelids. To move the face mesh according
9 to a set of sample points, control point positions must be computed that will
10 deform the surface appropriately. Using the same weighting functions described
11 above, we compute how the sample points move in response to the control points.
12 The result is a linear transformation: $\mathbf{p} = \mathbf{A}\mathbf{q}$. Therefore if at time t we want to
13 achieve the sample positions \mathbf{p}_t , we can use the control positions $\mathbf{q}_t = \mathbf{A}^{-1}\mathbf{p}_t$.
14 However, the matrix \mathbf{A} can be ill-conditioned, so to avoid the undesirable surface
15 shapes that are caused by very large control point motions we compute \mathbf{A}^{-1} using
16 the SVD (Singular Value Decomposition) and clamp the singular values of \mathbf{A}^{-1} at a
17 limit M . In the illustrated example, $M = 1.5$ was used. A standard reference that
18 discusses SVD is Golub and Van Loan, *Matrix Computations*, 3rd edition, Johns
19 Hopkins press, 1996.

20 21 **Eye and Head Movement**

22 In order to give the face a more lifelike appearance, procedurally generated
23 motion is added to the eyes and separately captured rigid-body motion to the head
24 as a whole. The eyeballs are rotated according to a random sequence of fixation
25 directions, moving smoothly from one to the next. The eyelids are animated by

1 rotating the vertices that define them about an axis through the center of the
2 eyeball, using weights defined on the eyelid mesh to ensure smooth deformations.

3 The rigid-body motion of the head is captured from the physical motion of
4 a person's head by filming that motion while the person is wearing a hat marked
5 with special machine-recognizable targets (the hat is patterned closely on the one
6 used by Marschner et al., *Image-based BRDF measurement including human skin*,
7 *Rendering Techniques '99* (Proceedings of the Eurographics Workshop on
8 *Rendering*), pps. 131-144, June 1998. By tracking these targets in the video
9 sequence, the rigid motion of the head is computed, which is then applied to the
10 head model for rendering. This setup, which requires simply a video camera,
11 provides a convenient way to author head motion by demonstrating the desired
12 actions.

13 14 **Exemplary System and Methods for Modeling Reflectance**

15 Rendering a realistic image of a face requires not just accurate geometry,
16 but also accurate computation of light reflection from the skin. In the illustrated
17 example, a physically-based Monte Carlo ray tracer was used to render the face.
18 Exemplary techniques are described in Cook et al., *Distribution Ray Tracing*,
19 *Computer Graphics (SIGGRAPH '84 Proceedings)*, pps. 165-174, July 1984 and
20 Shirley et al., *Monte Carlo techniques for direct lighting calculations*,
21 *Transactions on Graphics*, 15(1):1-36, 1996. Doing so allows for the use of
22 arbitrary BRDFs (bi-directional reflectance distribution functions) to correctly
23 simulate the appearance of the skin, which is not well approximated by simple
24 shading models. In addition, extended light sources are used, which, in rendering
25 as in portrait photography, are needed to achieve a pleasing image. Two important

1 *reflectance functions*, Computer Graphics (SIGGRAPH '97 Proceedings), pps.
2 117-126, August 1997.

3 4 **Constructing the Albedo Map**

5 In the illustrated and described embodiment, the albedo map, which must
6 describe the spatially varying reflectance due to diffuse reflection, was measured
7 using a sequence of digital photographs of the face taken under controlled
8 illumination.

9 Fig. 11 shows an exemplary system that was utilized to capture the digital
10 photographs or images. In the illustrated system, a digital camera 1100 is
11 provided and includes multiple light sources, exemplary ones of which are shown
12 at 1102, 1104. Polarizing filters in the form of perpendicular polarizers 1106,
13 1108, and 1110 are provided and cover the light sources and the camera lens so
14 that the specular reflections are suppressed, thereby leaving only the diffuse
15 component in the images. In the example, a subject wears a hat 1112 printed with
16 machine-recognizable targets to track head pose. Camera 1100 stays stationary
17 while the subject rotates. The only illumination comes from the light sources
18 1102, 1104 at measured locations near the camera. A black backdrop is used to
19 reduce indirect reflections from spilled light.

20 Since the camera and light source locations are known, standard ray tracing
21 techniques can be used to compute the surface normal, the irradiance, the viewing
22 direction, and the corresponding coordinates in texture space for each pixel in each
23 image. Under the assumption that ideal Lambertian reflection is being observed,
24 the Lambertian reflectance can be computed for a particular point in texture space
25 from this information. This computation is repeated for every pixel in one

second because the geometry provides visual detail—so this editing has relatively little effect on the appearance of the final renderings.

Fig. 16 shows several different aspects of the face model, using still frames from the accompanying video. In the first row, the face is shown from several angles to demonstrate that the albedo map and measured BRDF realistically capture the distinctive appearance of the skin and its color variation over the entire face, viewed from any angle. The second row shows the effects of rim and side lighting, including strong specular reflections at grazing angles. Note that the light source has the same intensity and is at the same distance from the face for all three images in this row. The directional variation in the reflectance leads to the familiar lighting effects seen in the renderings. In the third row, expression deformations are applied to the face to demonstrate that the face still looks natural under normal expression movement.

Fig. 17 is a flow diagram that describes steps in a method for creating an albedo map in accordance with the described embodiment. The method can be implemented in any suitable hardware, software, firmware or combination thereof. In the described embodiment, the method is implemented in software in connection with a system such as the one shown and described in Fig. 11.

Step 1700 provides one or more polarized light sources that can be used to illuminate a subject. Exemplary light sources are described above. In the described embodiment, the light sources are selected so that the specular component of the subject's facial reflectance is suppressed or eliminated. Step 1702 illuminates the subject's face with the light sources. Step 1704 rotates the subject while a series of digital photographs or images are taken. Step 1706 computes surface normals, irradiance, viewing direction and coordinates in texture

space for each pixel in the texture map. The computations can be done using known algorithms. Step 1708 computes the Lambertian reflectance for a particular pixel in the texture space for the image. This provides an albedo for the pixel. Step 1710 determines whether there are any additional pixels in the albedo map. If there are, step 1712 gets the next pixel and returns to step 1708. If there are no additional pixels in the albedo map, step 1714 ascertains whether there are any additional digital images. If there are additional digital images, step 1716 gets the next digital image and returns to step 1706. If there are no additional digital images, then step 1718 computes a weighted average of the individual albedo maps for each image to create a single albedo map for the entire face. One specific example of how this weighted average processing takes place is given above and described in Marschner, *Inverse Rendering for Computer Graphics*, PhD thesis, Cornell University, August 1998.

Fig. 18 is a flow diagram that describes steps in a method for computing an albedo for a single pixel. This method can be implemented in any suitable hardware, software, firmware or combination thereof. In the described embodiment, the method is implemented in software. Step 1800 determines, for a given pixel, whether the pixel is fully visible. If the pixel is not fully visible, then an albedo for the pixel is not computed (step 1804). If the pixel is fully visible, step 1802 determines whether the pixel is fully illuminated by at least one light source. If the pixel is not fully illuminated by at least one light source, then an albedo for the pixel is not computed (step 1804). If the pixel is fully illuminated by at least one light source, then step 1806 determines whether the pixel is partially illuminated by any light source. If so, then an albedo is not computed for the pixel. If the pixel is not partially illuminated by any light source, then step

1 1808 computes an albedo and a weight for the pixel. The weights are later used in
2 averaging together individual maps. Hence, as discussed above, albedos are
3 computed only for pixels that are fully visible, fully illuminated by at least one
4 light source, and not partially illuminated by any light source. This ensures that
5 partially occluded pixels and pixels that are in full-shadow or penumbra are not
6 used.

7 8 **Conclusion**

9 The embodiments described above provide systems and methods that
10 address the challenge of modeling and rendering faces to the high standard of
11 realism that must be met before an image as familiar as a human face can appear
12 believable. The philosophy of the approach is to use measurements whenever
13 possible so that the face model actually resembles a real face. The geometry of the
14 face is represented by a displacement-mapped subdivision surface that has
15 consistent connectivity and correspondence across different faces. The reflectance
16 comes from previous BRDF measurements of human skin together with new
17 measurements that combine several views into a single illumination-corrected
18 texture map for diffuse reflectance. The motion comes from previously described
19 motion capture technique and is applied to the face model using an improved
20 deformation method that produces motions suitable for shaded surfaces. The
21 realism of the renderings is greatly enhanced by using the geometry, motion, and
22 reflectance of real faces in a physically-based renderer.

23 Although the invention has been described in language specific to structural
24 features and/or methodological steps, it is to be understood that the invention
25 defined in the appended claims is not necessarily limited to the specific features or

steps described. Rather, the specific features and steps are disclosed as preferred forms of implementing the claimed invention.

forms of implementing the claimed invention.